## The NCBI Genomic Discovery System

Using Integrated Data from Multiple Genomes at NCBI to Accelerate Biological Discovery

Eric W. Sayers, Ph.D.

National Center for Biotechnology Information
Bethesda, MD USA

HGM2004 5 April 2004 17.00 – 19.00

# This Evening's Menu...

Gene of the Day: thyroid peroxidase

- Entrez Gene and RefSeq
- Exploring evidence for gene annotations
- Beyond RefSeq
  - What if my species isn't in RefSeq?
- Sequence Polymorphisms
  - Finding SNPs and phenotypes
- NCBI Map Viewer
- View and retrieve annotations and associated data
- Finding Gene Homologs
  - Find homologs using Entrez, Map Viewer, and BLAST
- Example Problems
  - Searching for markers and exploring disease alleles

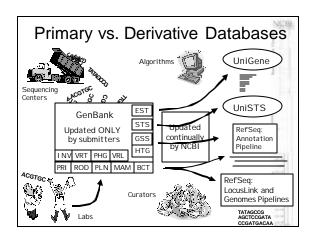
# The National Institutes of Health Bethesda, MD

# The National Center for Biotechnology Information • Created as a part of NLM in 1988 - Establish public databases - Perform research in computational biology - Develop software tools for sequence analysis - Disseminate biomedical information

# Entrez Gene and RefSeq

- Entrez Gene is a new database of gene loci
- Entrez Gene includes all organisms in LocusLink plus all other organisms in RefSeq
- Entrez Gene is an Entrez database with full searching and linking capabilities
- Entrez Gene will ultimately replace LocusLink

# Gene Record for Human TPO [In: IPO dayords perosidase. [floro apprent] [Completed 7/173. Lecu tag. Biology 2015 [September 2015] [Completed 7/173. Lecu tag. Biology 2015] [Complete 7/173. Lecu tag. Biology 2015] [Completed 7/173] [Comple



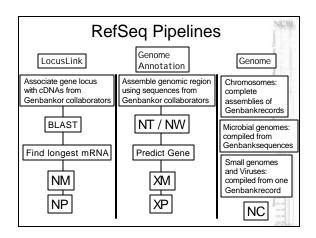
# RefSeq Benefits

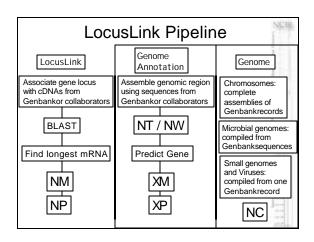
- Non-redundant
- · Explicitly linked nucleotide and protein sequences
- Updates to reflect current sequence data and biology
- · Data validation

XR 123456

- Format consistency
- Distinct accession series
- Stewardship by NCBI staff and collaborators

### RefSeq Accession Numbers Curated Records Complete genomic sequence (chromosome) Incomplete genomic sequence NC 123455 NG 123456 NM\_123456 mRNA NP 123456 Protein derived from NM NR\_123456 Non-coding RNA Model Records NT\_123456 Assembly of BAC data NW\_123456 Assembly of WGS data NZ\_ABCD12345678 Collection of WGS data XM\_123456 mRNA Protein derived from XM XP\_123456 Protein derived from NZ ZP\_123456 Non-coding RNA



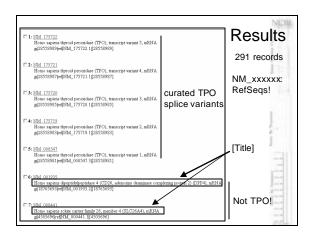


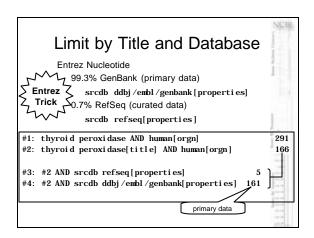
Curated Eu	ıka	ryotic RefSeqs	NCW.
Mus musculus	MM	Mus musculus	NΡ
Rattus norvegicus		Rattus norvegicus	-
Anopheles gambiae str. PEST		Anopheles gambiae str. PEST	. 34
Magnaporthe grisea 70-15		Magnaporthe grisea 70-15	
	NC	Eremothecium gossypii	11/1/17
	NC	Plasmodium falciparum 3D7	
Arabidopsis thaliana		Arabidopsis thaliana	
Oryza sativa (japonica cultivar-		Oryza sativa (japonica cultivar-group	0)
	NC	Schizosaccharomyces pombe	5 - 11
	NC	Saccharomyces cerevisiae	
	NC	Leishmania major	4
Trypanosoma brucei		Trypanosoma brucei Encephalitozoon cuniculi	3311
Caenorhabditis elegans	NC	Caenorhabditis elegans	11
		Drosophila melanogaster	127
Drosophila melanogaster Danio rerio		Danio rerio	133
Xenopus tropicalis		Xenopus tropicalis	2-4
Gallus gallus		Gallus gallus	10.3
Homo sapiens		Homo sapiens	-13
Bos taurus		Bos taurus	1311

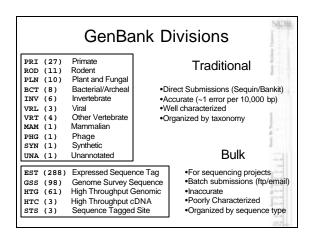
# Finding Primary Sequences

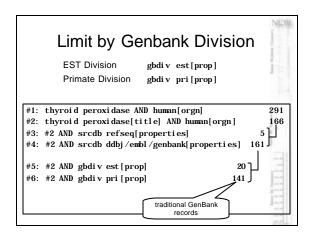
- Search Entrez Nucleotide
  - 99.3% GenBank (primary data)
  - 0.7% RefSeq (curated data)
- Typical starting query:
  - thyroid peroxidase AND human[orgn]
- Useful search terms in [Properties]:
  - srcdb : Source database (srcdb genbank)
  - gbdiv: GenBank division (gbdiv est)
  - biomol : Biomolecule type (biomol mrna)

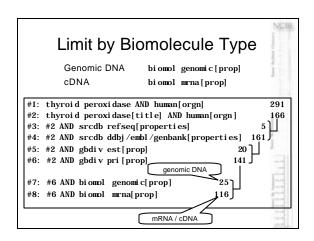
# Global Entrez Search Entrez, The Life Sciences Search Engine Entrez, The Life Sciences Search Engine Search across databases Search across databases Search across databases Fundament birmelical literature citations and encircles Search across databases Fundament of the search science of the science of the

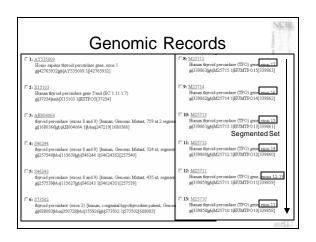


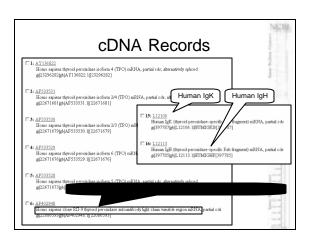


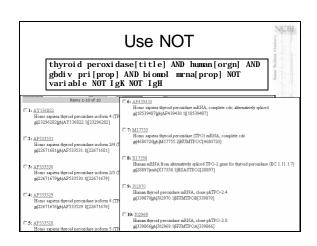


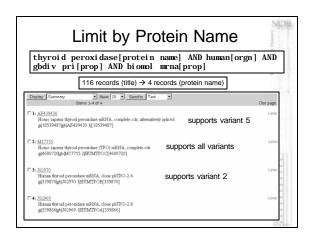


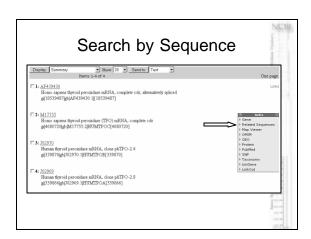


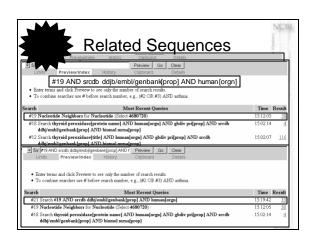


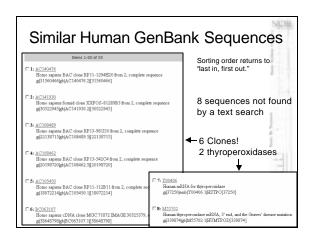


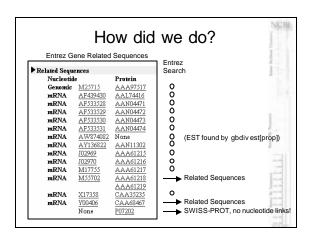


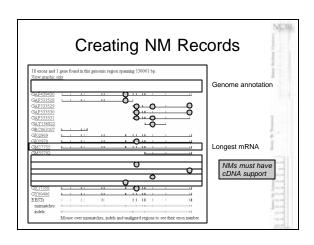


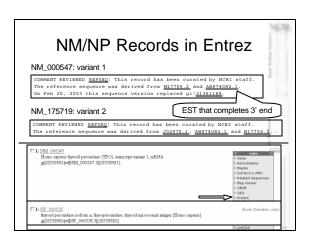


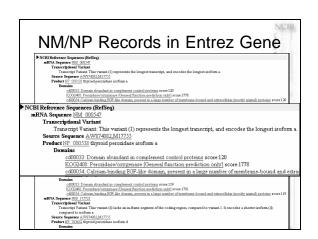


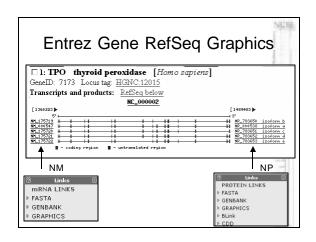


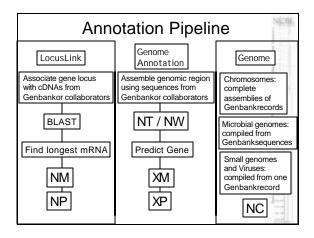


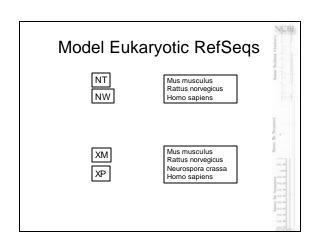


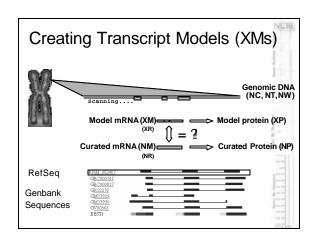


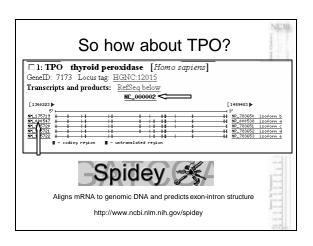


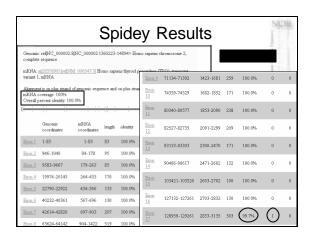


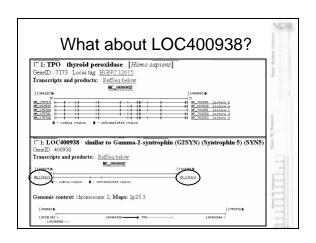


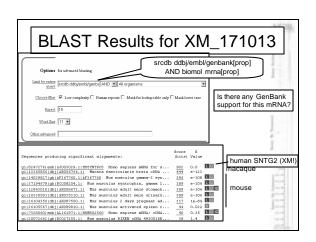


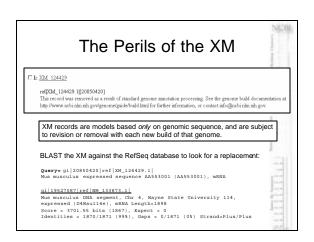


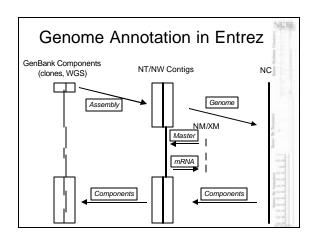


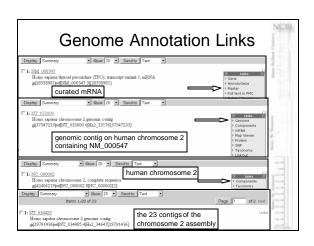


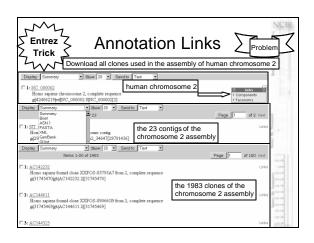


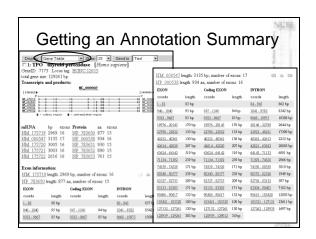


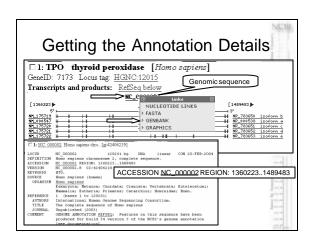


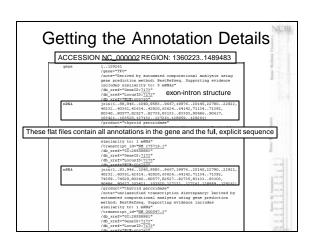


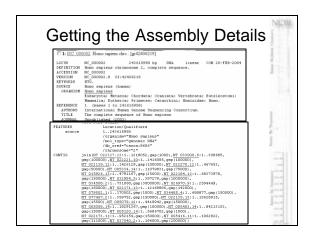


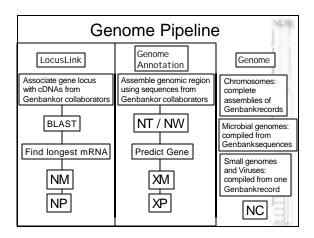


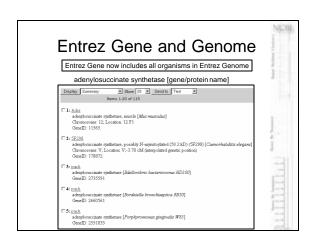


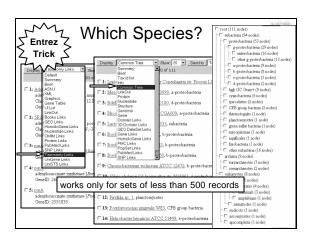


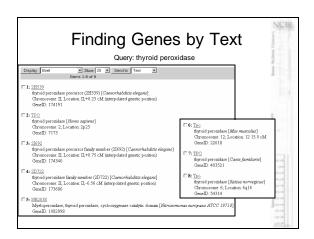




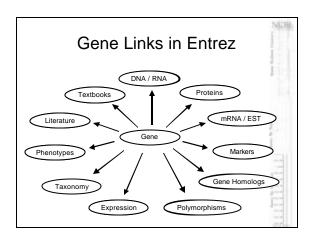


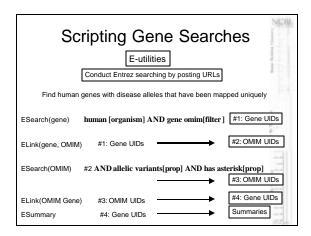






Searching Entrez Gene
Gene symbol: human thyroid peroxidase (TPO)
tpo [sym] AND human [organism]
Protein name: topoisomerase genes from Archaea
topoisomerase[gene/protein name] AND archaea [organism]
Chromosome and Links: genes on human chromosome 2 with OMIM links 2 [chromosome] AND gene omim [filter] AND human [organism]
RefSeq status and variants: Reviewed RefSeqs with transcript variants srcdb refseq reviewed[prop] AND has transcript variants[prop]
Disease and Gene Ontology: Membrane proteins linked to cancer integral to plasma membrane[gene ontology] AND cancer [dis]





# Gene, LocusLink, and the Future

- · Gene will eventually replace LocusLink
- Many LocusLink functions are now in the Entrez Links
   menu
- As much as possible, LocusLink IDs are equal to the corresponding Gene UID
- Currently there is no FTP site for Entrez Gene
- Therefore, we will maintain the LL\_tmpl file and associated files until Entrez Gene has FTP files

# Beyond RefSeq

If your organism does not have RefSeqs...

• UniGene : gene-based clusters of cDNAs and ESTs

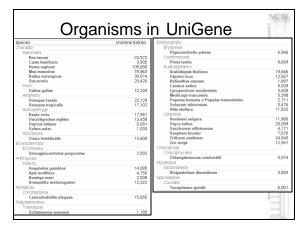
• UniSTS : STS markers

• WGS sequences in Entrez Nucleotide (wgs[prop])

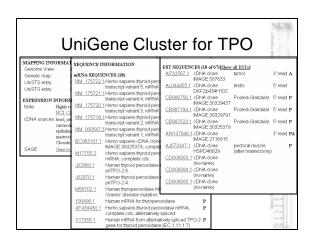
Trace Archive

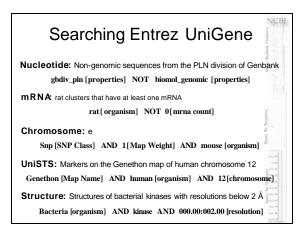
# UniGene: Clustering Expressed Sequences

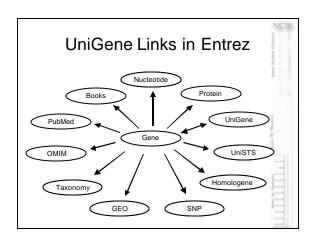
- Records are clusters of mRNAs and ESTs that ideally represent single genes
- Records are created automatically by a modified BLAST algorithm
- UniGene provides a means to identify an EST or unannotated mRNA

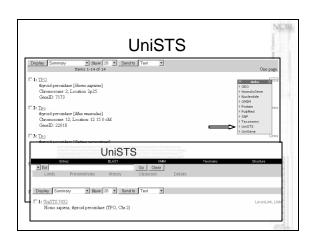


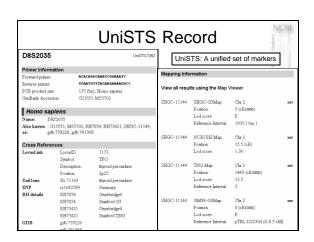
F1: TPO thyroid peroxidase [Homo supters]   GeneD 7173 Locus lag HSRC12012   Transcripts and products: McSea below [Lineary [Li	b GEO  h MomoloGene  Nucleotide  Nucleotide
Genomic contexts: chromosome 2; Maps: 2p25	► PubMed  ► SMP  ► Taxonomy  ► UniSTS  ► AceView  ► HGNC  ► KEGG
[195013] 10000000 — 10000000 — 10000000 — 10000000 — 100000000	► LocusID ► Map Viewer ► UniGene
▼ for M17755 Go Clear Limits PreviewIndex History Cipboard Details by Er  Diselew Summery ▼ Shower 20 ▼ Sort ▼ Seed to Text ▼	ntrez search

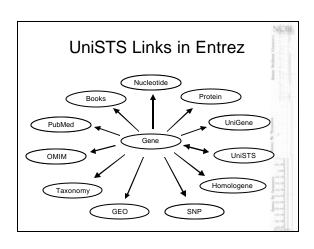


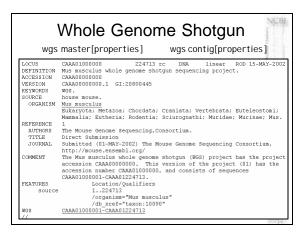


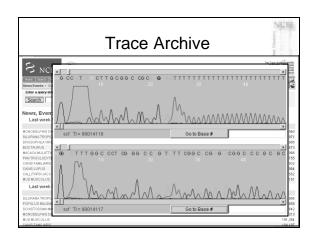


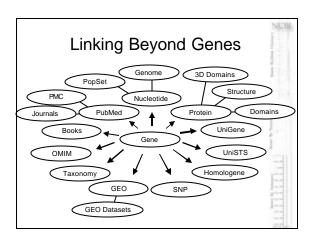


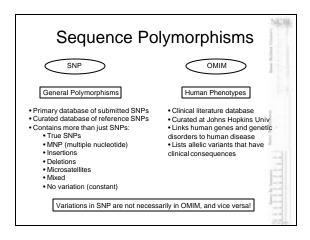


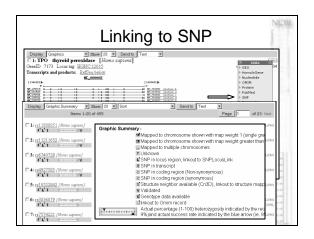


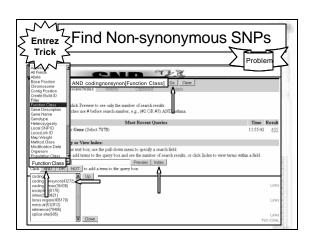


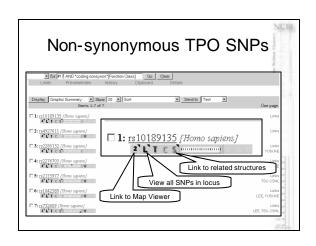


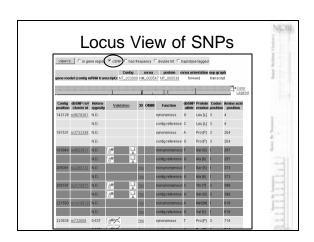


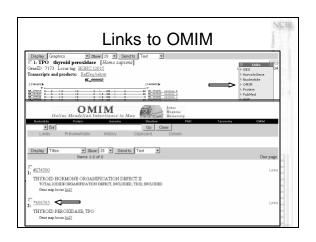


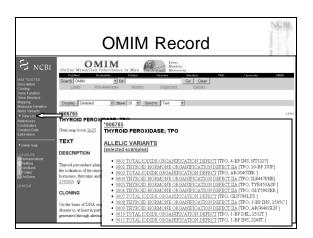


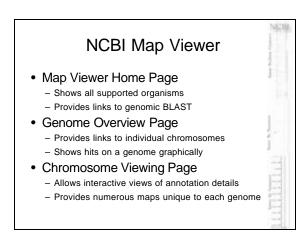


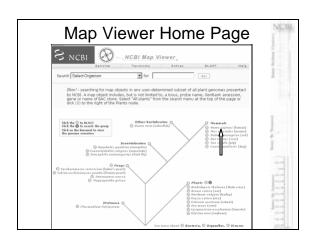


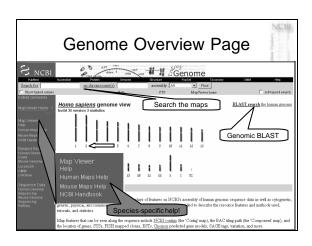


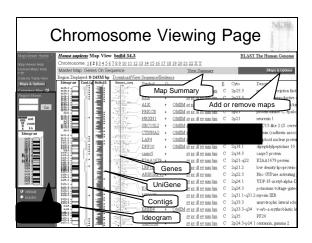


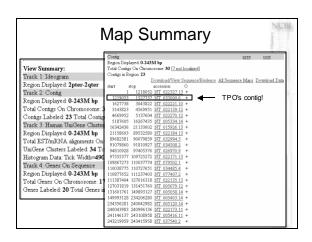


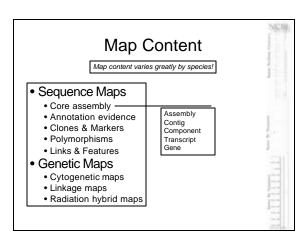


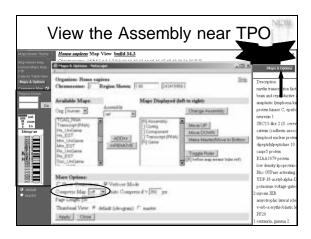




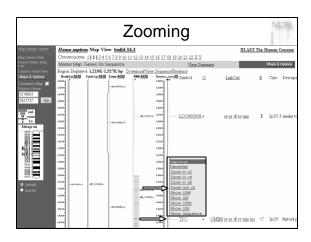


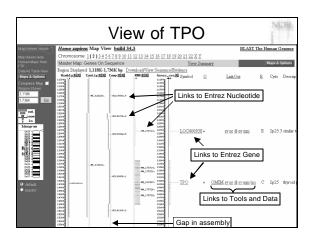


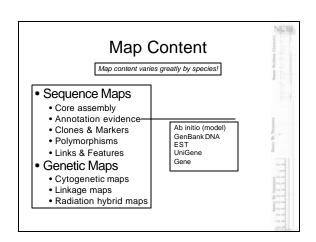


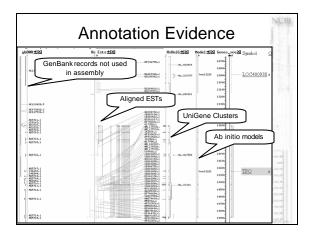


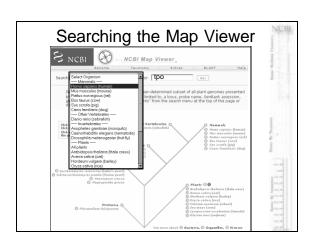
		As	ser	nbly	/ of	Cł	٦r	. 2		18	14.
Map Viewer Home Map Viewer Help Human Maps Help FTP	Homo sapie Chromosome Master Map:	: 1[2]34 Genes On:	5678910: Sequence	11 12 13 14 15		View	Y Summa	BLAST		uman (	
Page 12 Colored Topics of Colo	Region Display  Rand by 2120  120000  120000  120000  120000  120000  120000  120000  120000		527K bp T Core 236 Core 336 Co	M_17719	SOURCE OF	indence  Symbol  LOC4009		Lank Out  or of 40 er com	E	Cyto 2p25.3	Descri

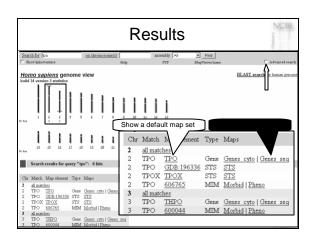


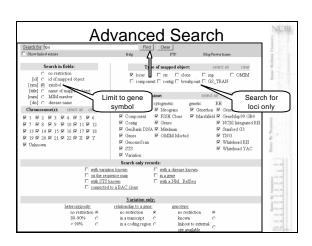


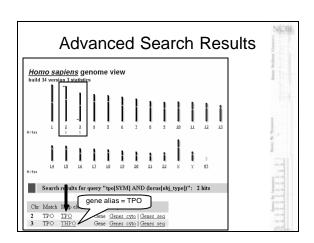


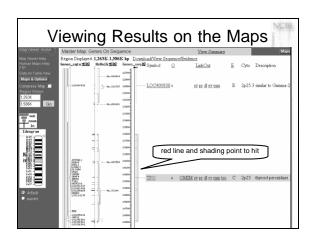


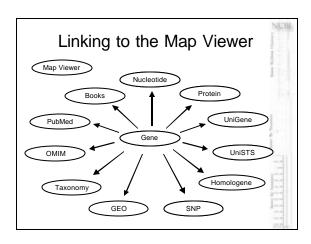


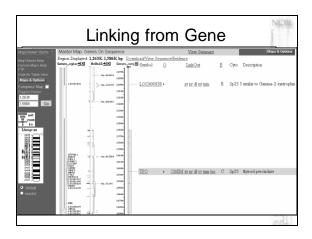


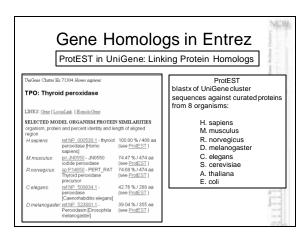


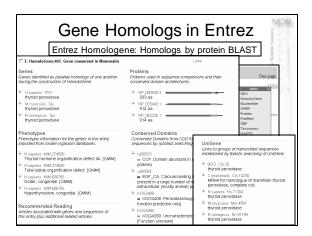


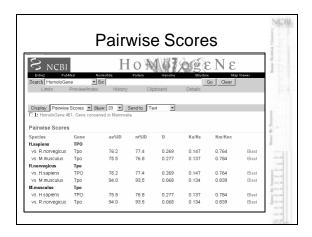


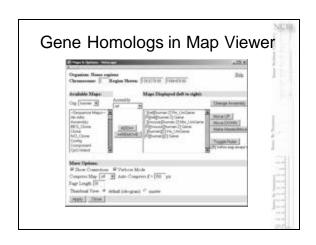


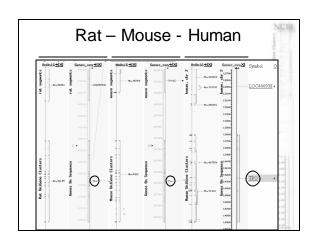


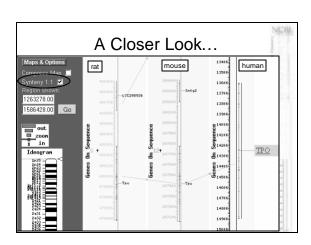


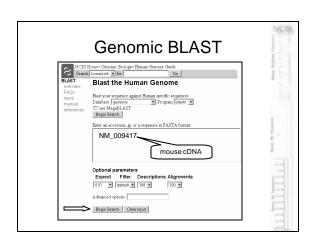


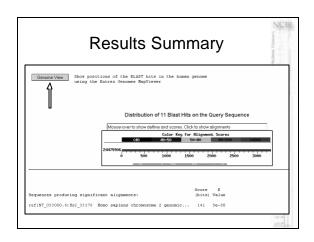


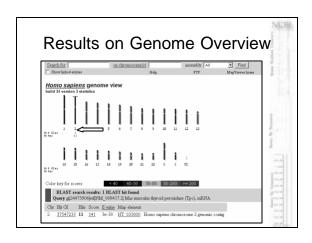


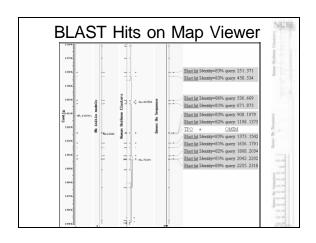


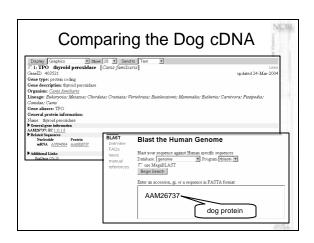


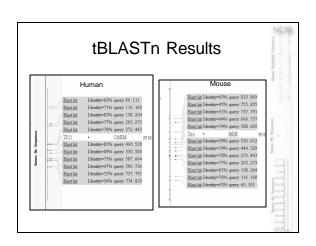


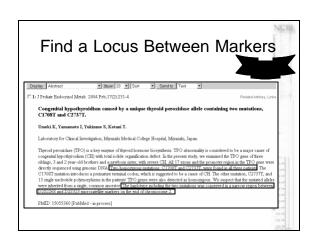


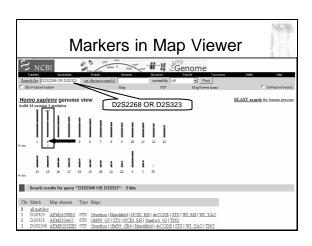


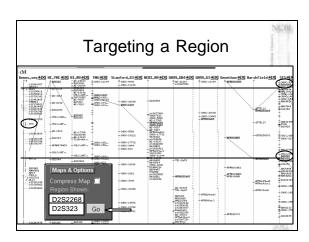


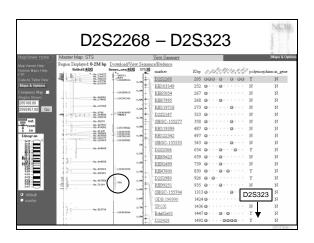


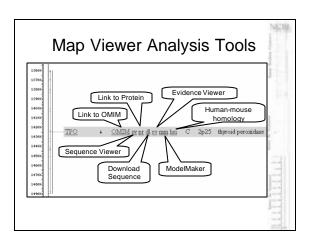


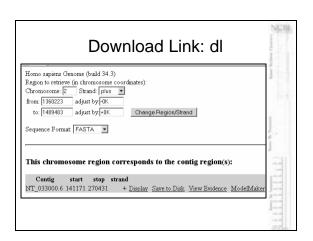


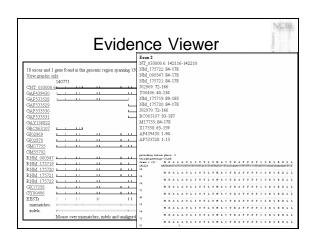


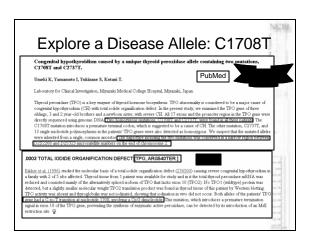


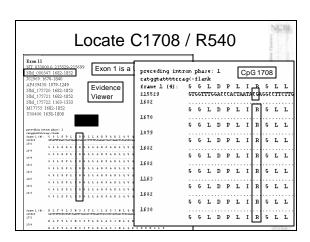


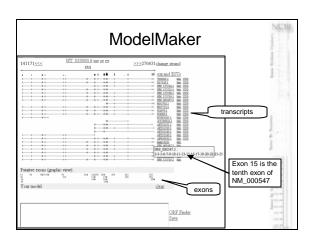


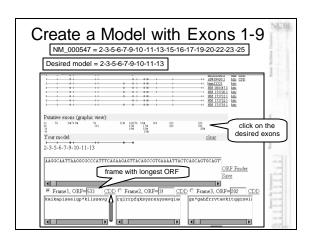


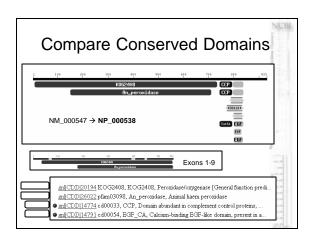


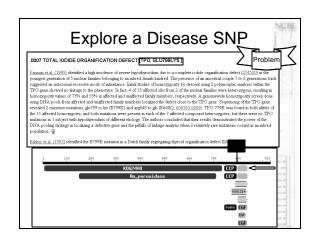


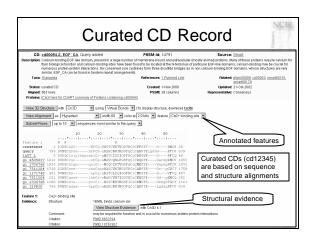


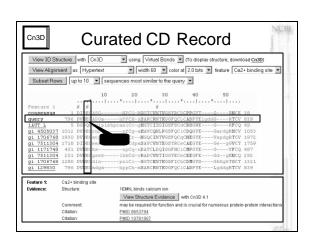












For Mo	re Information
• Eric Sayers • General Help • BLAST	sayers@ncbi.nlm.nih.gov info@ncbi.nlm.nih.gov blast-help@ncbi.nlm.nih.gov
The NCBI Hair Follow the link	from the NCBI Home Page
Co	ome See Us!
	inc 000 03:
	Booth #12